

PARTIALLY OBSERVABLE RISK-SENSITIVE MARKOV DECISION PROCESSES

NICOLE BÄUERLE* AND ULRICH RIEDER†

ABSTRACT. We consider the problem of minimizing a certainty equivalent of the total or discounted cost over a finite and an infinite time horizon which is generated by a Partially Observable Markov Decision Process (POMDP). The certainty equivalent is defined by $U^{-1}(\mathbb{E} U(Y))$ where U is an increasing function. In contrast to a risk-neutral decision maker, this optimization criterion takes the variability of the cost into account. It contains as a special case the classical risk-sensitive optimization criterion with an exponential utility. We show that this optimization problem can be solved by embedding the problem into a completely observable Markov Decision Process with extended state space and give conditions under which an optimal policy exists. The state space has to be extended by the joint conditional distribution of current unobserved state and accumulated cost. In case of an exponential utility, the problem simplifies considerably and we rediscover what in previous literature has been named *information state*. However, since we do not use any change of measure techniques here, our approach is simpler. A simple example, namely a risk-sensitive Bayesian house selling problem is considered to illustrate our results.

KEY WORDS: Partially Observable Markov Decision Problem, Certainty Equivalent, Exponential Utility, Updating Operator, Value Iteration.

1. INTRODUCTION

In this work we consider Partially Observable Markov Decision Processes (POMDP) under a general risk-sensitive optimization criterion for problems with finite and infinite time horizon. This is a continuation of our research published in [2]. More precisely our aim is to minimize the certainty equivalent of the accumulated total cost of a POMDP. In case of an infinite time horizon, costs have to be discounted. The certainty equivalent of a random variable is defined by $U^{-1}(\mathbb{E} U(X))$ where U is an increasing function. If $U(x) = x$ we obtain as a special case the classical risk-neutral decision maker. The case $U(x) = \frac{1}{\gamma} e^{\gamma x}$ is often referred to as 'risk-sensitive', however the risk-sensitivity is here only expressed in a special way through the risk-sensitivity parameter $\gamma \neq 0$. More general, the certainty equivalent may be written (assuming enough regularity of U) as

$$U^{-1}\left(\mathbb{E}[U(X)]\right) \approx \mathbb{E} X - \frac{1}{2} l_U(\mathbb{E} X) \text{Var}[X] \quad (1.1)$$

where

$$l_U(x) = -\frac{U''(x)}{U'(x)}$$

is the *Arrow-Pratt* function of absolute risk aversion. In case of an exponential utility, this absolute risk aversion is constant (for a discussion see [5]). If U is concave, the variance is subtracted and the decision maker is risk seeking in case cost is minimized, if U is convex, then the variance is added and the decision maker is risk averse.

In case of complete observation it has been shown in [2] that this problem can be recast in the theory of Markov Decision Processes (MDP) by enlarging the state space with the total discounted cost that has been incurred so far. Numerical solution procedures via linear programming of these completely observable general risk-sensitive Markov Decision Processes can be found in [10]. The average cost version of this problem is treated in [7] and for an application

in insurance see [3]. Now we assume that only one of two components of a controlled Markov process can be observed. However, also the cost may depend on both components which leads to the situation that the cost incurred so far is an unobservable quantity. It is well-known that in case of a risk-neutral decision maker, the partially observable problem can be solved by a completely observable MDP when we enlarge the state space by the conditional distribution of the unobservable state, given the observable history of the process (see e.g. [1] chapter 5, [12] chapter 4 or [13] chapter 7). As far as the risk-sensitive problem is concerned we proceed in a similar way. This time however, the corresponding problem with complete observation possesses already an enlarged state consisting of the process state and the total discounted cost so far. Thus, to cope with the partially observable model we construct a Markov Decision Process where the state consists of the observable part of the state and the joint conditional distribution of the unobservable part of the state and unobservable total cost so far, given the observable history of the process.

Early papers [17, 20] provided rigorous mathematical treatment of POMDPs with Borel state and action spaces. These references already present the solution procedure via the enlargement of the state space and the reduction to an ordinary Markov Decision Process. For a detailed discussion of the theory in the classical risk-neutral setting and for several applications see [1] chapter 5 and [12] chapter 4. Risk-sensitive Markov Decision processes with the exponential utility have been discussed intensively since the seminal paper of [14]. For further references we refer the reader to [2]. Recent applications of this criterion in a wide range of portfolio optimization problems can be found in [8]. Papers which combine the exponential utility with POMDPs are among others [15, 11, 9, 18, 6]. In all these papers a control model formulation has been used, where the true, unobservable and controlled state process is a Markov process (under Markovian policies) and observations are obtained by perturbed signals of this process. A change of measure technique is used to obtain independent signals. In order to apply MDP theory, the state space has been enlarged by a quantity that has been called an 'information vector'. In the present paper we use a more general model formulation where both parts (observable and unobservable state) are jointly Markovian and can be controlled jointly. This setting also covers the Bayesian case where the unknown state part is simply an unknown parameter. Also note that our optimization criterion is not restricted to the exponential utility and we do not need a change of measure technique to derive our filter. Moreover, the general approach implies a very natural interpretation for the 'information vector' in the exponential utility case. Besides [15] all the previously mentioned papers focus on the risk-sensitive average criterion by using the vanishing discount approach, i.e., by looking at the β -discounted problem and by letting β go to 1. In [6] a finite state and action space is considered and emphasis is laid on numerical aspects of the problem. A discrete-time linear quadratic risk-sensitive stochastic control problem with incomplete state information is solved in [19].

Our paper is organized as follows: In the next section we introduce the underlying POMDP and define general history-dependent (deterministic) policies for this model. In section 3 we consider the finite horizon general risk-sensitive problem and introduce continuity and compactness assumptions which will guarantee the existence of optimal policies. Then the problem is embedded into a suitably defined Markov Decision Process where the state space contains among others a joint conditional distribution of the unobservable state and total accumulated cost so far, given the observed process. An updating-operator is defined to create a forward iteration of this joint conditional distribution. The main theorem of this section (Theorem 3.3) states the validity of the embedding procedure and the existence of optimal policies. Section 4 contains some important special cases. Among them the situation where the cost function does not depend on the unobservable state in which case the updating operator simplifies to the updating operator for classical risk-neutral POMDPs. In case the exponential utility function is used, we rediscover some results of the previous literature. We also consider the case of a power utility where we only get a slight simplification. In Section 5 we consider a simple risk-sensitive Bayesian house selling problem. We prove the existence of so-called 'reservation levels' which

can be seen as thresholds for the acceptance of an offer. These reservation levels depend only on the conditional distribution. In the last section we consider the problem with infinite time horizon and distinguish the case of a convex and a concave utility functions which require separate proofs due to different inequalities. The main theorems (Theorem 6.1, Theorem 6.2) show that the value function of the problem can be obtained from a fixed point equation and that an optimal policy exists which is not stationary but still can be generated by only one decision function.

2. GENERAL PARTIALLY OBSERVABLE RISK-SENSITIVE MARKOV DECISION PROCESSES

We suppose that a *partially observable Markov Decision Processes* is given which we introduce as follows: We denote this process by $(X_n, Y_n)_{n \in \mathbb{N}_0}$ and assume that the state space is $E_X \times E_Y$ where E_X and E_Y are Borel spaces, i.e., Borel subsets of some Polish spaces. The x -component will be the *observable* part, the y -component *cannot be observed* by the controller. Actions can be taken from a set A which is again a Borel space. The set $D \subset E_X \times A$ is a Borel subset of $E_X \times A$. By $D(x) := \{a \in A : (x, a) \in D\}$ we denote the feasible actions depending on the observable state part x . We assume that D contains the graph of a measurable mapping from E_X to A . There is a stochastic transition kernel Q from $D \times E_Y$ to $E_X \times E_Y$ which determines the distribution of the new state pair given the current state and action. So $Q(B|x, y, a)$ is the probability that the next state pair is in $B \in \mathcal{B}(E_X \times E_Y)$, given the current state is (x, y) and action $a \in D(x)$ is taken. In what follows we assume that the transition kernel Q has a measurable density q with respect to some σ -finite measures λ and ν , i.e.,

$$Q(B|x, y, a) = \int_B q(x', y'|x, y, a) \lambda(dx') \nu(dy'), \quad B \in \mathcal{B}(E_X \times E_Y).$$

For convenience we introduce the marginal transition kernel density by

$$q^X(x'|x, y, a) := \int_{E_Y} q(x', y'|x, y, a) \nu(dy').$$

We assume that the initial distribution Q_0 of Y_0 is known. Further we have a measurable one-stage cost function $c : D \times E_Y \rightarrow \mathbb{R}_+$. We assume in particular that the cost $c(x, y, a)$ also depends on the unknown state part y . Finally we have a discount factor $\beta \in (0, 1]$.

Next we introduce policies for the controller. Here it is important to consider the *set of observable histories* which are defined as follows:

$$\begin{aligned} H_0 &:= E_X \\ H_n &:= H_{n-1} \times A \times E_X. \end{aligned}$$

An element $h_n = (x_0, a_0, x_1, \dots, x_n) \in H_n$ denotes the observable history of the process up to time n .

Definition 2.1. a) A measurable mapping $g_n : H_n \rightarrow A$ with the property $g_n(h_n) \in D(x_n)$ for $h_n \in H_n$ is called a *decision rule* at stage n .

b) A sequence $\pi = (g_0, g_1, \dots)$ where g_n is a decision rule at stage n for all n , is called *policy*. We denote by Π the set of all policies.

3. FINITE HORIZON PROBLEMS

In this section we consider problems with finite time horizon N . For a fixed policy $\pi = (g_0, g_1, \dots) \in \Pi$ and fixed (observable) initial state $x \in E_X$, the initial distribution Q_0 together with the transition kernel Q define by a theorem of Ionescu Tulcea a probability measure \mathbb{P}_{xy}^π on $(E_X \times E_Y)^{N+1}$ endowed with the product σ -algebra. More precisely \mathbb{P}_{xy}^π is the probability measure under policy π given $X_0 = x$ and $Y_0 = y$. Later we also use the probability measure $\mathbb{P}_x^\pi(\cdot) := \int \mathbb{P}_{xy}^\pi(\cdot) Q_0(dy)$. For $\omega = (x_0, y_0, \dots, x_N, y_N) \in (E_X \times E_Y)^{N+1}$ we define the random variables X_n and Y_n in a canonical way by their projections

$$X_n(\omega) = x_n, \quad Y_n(\omega) = y_n.$$

If $\pi = (g_0, g_1, \dots) \in \Pi$ is a given policy, we define recursively

$$\begin{aligned} A_0 &:= g_0(X_0) \\ A_n &:= g_n(X_0, A_0, X_1, \dots, X_n), \end{aligned}$$

the sequence of actions which are chosen successively under policy π . We assume that the decision maker is risk averse and has a utility function $U : \mathbb{R}_+ \rightarrow \mathbb{R}$ which is continuous and strictly increasing. The optimization problem is defined as follows. For $\pi \in \Pi$ and $X_0 = x$ denote

$$J_{N\pi}(x) := \int_{E_Y} \mathbb{E}_{xy}^\pi \left[U \left(\sum_{k=0}^{N-1} \beta^k c(X_k, Y_k, A_k) \right) \right] Q_0(dy)$$

and

$$J_N(x) := \inf_{\pi \in \Pi} J_{N\pi}(x). \quad (3.1)$$

Note that in case $U(x) = x$ we end up with the usual risk neutral Partially Observable Markov Decision Process setup (see e.g. [1] chapter 5, [12] chapter 4). Here however, if U is strictly concave, then U is a utility function and $U^{-1}(J_N(x))$ represents a *certainty equivalent*. If U is concave, we can see from (1.1) that the decision maker is risk seeking and if U is convex, then the decision maker is risk averse.

In what follows we show how to solve these kind of problems by using an *embedding technique*. In order to later ensure the existence of integrals and optimal policies we make the following assumptions (A):

- (i) $U : [0, \infty) \rightarrow \mathbb{R}$ is continuous and strictly increasing,
- (ii) $D(x)$ is compact for all $x \in E_X$,
- (iii) $x \mapsto D(x)$ is upper semicontinuous, i.e. for all $x \in E_X$ it holds: If $x_n \rightarrow x$ and $a_n \in D(x_n)$ for all $n \in \mathbb{N}$, then (a_n) has an accumulation point in $D(x)$,
- (iv) $(x, y, a) \mapsto c(x, y, a)$ is continuous,
- (v) $(x, y, x', y', a) \mapsto q(x', y' | x, y, a)$ is continuous and bounded.
- (vi) c is bounded, i.e., there exist constants $0 < \underline{c} < \bar{c}$ with $\underline{c} \leq c(x, y, a) \leq \bar{c}$.

Remark 3.1. Note that these assumptions are quite strong, however include in particular the case when state and action spaces are finite. (A)(ii-v) also ensure the existence of optimal policies for risk-neutral POMDP.

In [2] we have solved problem (3.1) for the observable case by extending the state space to include the accumulated cost so far. Now in the unobservable model, the state y and the accumulated cost so far cannot be observed because it depends on y . Thus, we proceed as in risk-neutral POMDPs (see e.g. [17, 20]) and consider probability measures μ on $E_Y \times \mathbb{R}_+$:

$$\begin{aligned} \mu \in \mathbb{P}_b(E_Y \times \mathbb{R}_+) &:= \left\{ \mu \text{ is a probability measure on the } \sigma\text{-algebra } \mathcal{B}(E_Y \times \mathbb{R}_+) \text{ such} \right. \\ &\quad \left. \text{that there exists a constant } K = K(\mu) > 0 \text{ with } \mu(E_Y \times [0, K]) = 1 \right\}. \end{aligned}$$

μ plays the role of the conditional distribution on the larger state space of hidden state component and accumulated cost. The precise interpretation will be seen in Theorem 3.2. In order to solve the optimization problem, we need, as in the risk-neutral case, an updating procedure for the conditional distributions which generates the filter process. The following updating-operator $\Psi : E_X \times A \times E_X \times \mathbb{P}_b(E_Y \times \mathbb{R}_+) \times \mathbb{R}_+ \rightarrow \mathbb{P}_b(E_Y \times \mathbb{R}_+)$ will do the task:

$$\Psi(x, a, x', \mu, z)(B) := \frac{\int_{E_Y} \int_{\mathbb{R}_+} \left(\int_B q(x', y' | x, y, a) \nu(dy') \delta_{s+zc(x, y, a)}(ds') \right) \mu(dy, ds)}{\int_{E_Y} q^X(x' | x, y, a) \mu^Y(dy)} \quad (3.2)$$

where $B \in \mathcal{B}(E_Y \times \mathbb{R}_+)$ and $\mu^Y(dy) := \mu(dy, \mathbb{R}_+)$ is the Y -marginal distribution of μ . Later we will also need the S -marginal $\mu^S(ds) := \mu(E_Y, ds)$. We define the updating operator only when

the denominator is positive. For $n \in \mathbb{N}$, $h_n := (x_0, a_0, \dots, x_n)$ and $B \in \mathcal{B}(E_Y \times \mathbb{R}_+)$ define now a sequence of probability measures

$$\begin{aligned}\mu_0(B|h_0) &:= (Q_0 \otimes \delta_0)(B), \\ \mu_{n+1}(B|h_n, a, x') &:= \Psi(x_n, a, x', \mu_n(\cdot|h_n), \beta^n)(B).\end{aligned}\tag{3.3}$$

The next theorem shows that the sequence of probability measures (μ_n) has the intended interpretation. For this purpose define the r.v.

$$S_0 := 0, \quad S_n := \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k), \quad n \in \mathbb{N}.$$

We then obtain:

Theorem 3.2. *Suppose (μ_n) is given by the recursion (3.3). For $n \in \mathbb{N}_0$ and all $\pi \in \Pi$ it holds that*

$$\mathbb{P}_x^\pi((Y_n, S_n) \in B | X_0, A_0, \dots, X_n) = \mu_n(B | X_0, A_0, \dots, X_n) \quad \mathbb{P}_x^\pi - a.s., \quad \text{for } B \in \mathcal{B}(E_Y \times \mathbb{R}_+).$$

Proof. Recall that $\mathbb{P}_x^\pi(\cdot) := \int \mathbb{P}_{xy}^\pi(\cdot) Q_0(dy)$. We first show that

$$E_x^\pi \left[v(X_0, A_0, X_1, \dots, X_n, Y_n, S_n) \right] = E_x^\pi \left[v'(X_0, A_0, X_1, \dots, X_n) \right] \tag{3.4}$$

for all bounded and measurable $v : H_n \times E_Y \times \mathbb{R}_+ \rightarrow \mathbb{R}$ and

$$v'(h_n) := \int_{E_Y} \int_{\mathbb{R}_+} v(h_n, y_n, s_n) \mu_n(dy_n, ds_n | h_n).$$

We do this by induction. For $n = 0$ both sides reduce to $\int v(x, y, 0) Q_0(dy)$. Now suppose the statement is true for $n - 1$. We simply write g_n instead of $g_n(h_n)$. We obtain for the left-hand side with a given observable history h_{n-1} :

$$\begin{aligned}E_x^\pi \left[v(h_{n-1}, A_{n-1}, X_n, Y_n, S_n) \right] &= \int_{E_Y} \int_{\mathbb{R}_+} \mu_{n-1}(dy_{n-1}, ds_{n-1} | h_{n-1}) \\ &\quad \cdot \int_{E_Y} \int_{E_X} \nu(dy_n) \lambda(dx_n) q(x_n, y_n | x_{n-1}, y_{n-1}, g_{n-1}) \\ &\quad \cdot \int_{\mathbb{R}_+} \delta_{s_{n-1} + \beta^{n-1} c(x_{n-1}, y_{n-1}, g_{n-1})}(ds_n) v(h_{n-1}, g_{n-1}, x_n, y_n, s_n) \\ &= \int_{E_Y} \int_{\mathbb{R}_+} \mu_{n-1}(dy_{n-1}, ds_{n-1} | h_{n-1}) \int_{E_Y} \int_{E_X} \nu(dy_n) \lambda(dx_n) q(x_n, y_n | x_{n-1}, y_{n-1}, g_{n-1}) \\ &\quad \cdot v(h_{n-1}, g_{n-1}, x_n, y_n, s_{n-1} + \beta^{n-1} c(x_{n-1}, y_{n-1}, g_{n-1})).\end{aligned}$$

For the right-hand side we obtain (where we insert the recursion for μ_n in the third equation and use Fubini's theorem, so that the normalizing constant of μ_n cancels out):

$$\begin{aligned}
E_x^\pi \left[v'(h_{n-1}, A_{n-1}, X_n) \right] &= \int_{E_Y} \int_{\mathbb{R}_+} \mu_{n-1}(dy_{n-1}, ds_{n-1} | h_{n-1}) \\
&\quad \cdot \int_{E_X} \lambda(dx_n) q^X(x_n | x_{n-1}, y_{n-1}, g_{n-1}) v'(h_{n-1}, g_{n-1}, x_n) \\
&= \int_{E_Y} \mu_{n-1}^Y(dy_{n-1} | h_{n-1}) \int_{E_X} \lambda(dx_n) q^X(x_n | x_{n-1}, y_{n-1}, g_{n-1}) \\
&\quad \cdot \int_{E_Y} \int_{\mathbb{R}_+} \mu_n(dy_n, ds_n | h_n) v(h_{n-1}, g_{n-1}, x_n, y_n, s_n) \\
&= \int_{E_Y} \int_{E_X} \nu(dy_n) \lambda(dx_n) \int_{E_Y} \int_{\mathbb{R}_+} \mu_{n-1}(dy_{n-1}, ds_{n-1} | h_{n-1}) q(x_n, y_n | x_{n-1}, y_{n-1}, g_{n-1}) \\
&\quad \cdot \int_{\mathbb{R}_+} \delta_{s_{n-1} + \beta^{n-1} c(x_{n-1}, y_{n-1}, g_{n-1})}(ds_n) v(h_{n-1}, g_{n-1}, x_n, y_n, s_n) \\
&= \int_{E_Y} \int_{\mathbb{R}_+} \mu_{n-1}(dy_{n-1}, ds_{n-1} | h_{n-1}) \int_{E_Y} \int_{E_X} \nu(dy_n) \lambda(dx_n) q(x_n, y_n | x_{n-1}, y_{n-1}, g_{n-1}) \\
&\quad \cdot v(h_{n-1}, g_{n-1}, x_n, y_n, s_{n-1} + \beta^{n-1} c(x_{n-1}, y_{n-1}, g_{n-1})).
\end{aligned}$$

Thus equation (3.4) is proved. It implies in particular for $v = 1_{B \times C}$ with $B \in \mathcal{B}(E_Y \times \mathbb{R}_+)$ and $C \subset E_X \times A \times \dots \times E_X$ a measurable set of histories until time n that

$$\mathbb{P}_x^\pi((Y_n, S_n) \in B, (X_0, A_0, \dots, X_n) \in C) = \mathbb{E}_x^\pi[\mu_n(B | X_0, A_0, \dots, X_n) 1_C((X_0, A_0, \dots, X_n))].$$

This in turn yields by definition that $\mu_n(B | X_0, A_0, \dots, X_n)$ is a conditional \mathbb{P}_x^π -distribution of (Y_n, S_n) given the history (X_0, A_0, \dots, X_n) . \square

Now we turn again to the optimization problem (3.1). Motivated by the previous result we define for $x \in E_X$, $\mu \in \mathbb{P}_b(E_Y \times \mathbb{R}_+)$, $z \in (0, 1]$ and $n = 1, \dots, N$:

$$V_{n\pi}(x, \mu, z) := \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}_{xy}^\pi \left[U \left(s + z \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k) \right) \right] \mu(dy, ds) \quad (3.5)$$

$$V_n(x, \mu, z) := \inf_{\pi \in \Pi} V_{n\pi}(x, \mu, z). \quad (3.6)$$

Obviously we have that $J_N(x) = V_N(x, Q_0 \otimes \delta_0, 1)$ where δ_x is the Dirac-measure at the point $x \in \mathbb{R}$. However, problem (3.6) can be solved with the general theory of POMDP and [2] by defining a suitable MDP. For this purpose let us define for a probability measure $\mu \in \mathbb{P}(E_Y)$

$$Q^X(B | x, \mu, a) := \int_B \int_{E_Y} q^X(x' | x, y, a) \mu(dy) \lambda(dx'), \quad B \in \mathcal{B}(E_X)$$

We consider a Markov Decision Process with state space $E := E_X \times \mathbb{P}_b(E_Y \times \mathbb{R}_+) \times (0, 1]$, action space A and admissible actions given by the set D . The one-stage cost is zero and the terminal cost function is $V_0(x, \mu, z) := \int \int U(s) \mu(dy, ds)$. Note that for all $\mu \in \mathbb{P}_b(E_Y \times \mathbb{R}_+)$ the expectation is well-defined since the support of μ in the s -component is a compact set. The transition law is given by $\tilde{Q}(\cdot | x, \mu, z, a)$ which is for $(x, \mu, z, a) \in E \times A$, $a \in D(x)$ and a measurable subset $B \subset E$ defined by

$$\tilde{Q}(B | x, \mu, z, a) := \int_{E_X} 1_B((x', \Psi(x, a, x', \mu, z), \beta z)) Q^X(dx' | x, \mu^Y, a).$$

Note that \tilde{Q} is again a transition kernel. Decision rules in the MDP setting are given by measurable mappings $f : E \rightarrow A$ such that $f(x, \mu, z) \in D(x)$. We denote by F the set of decision rules and by Π^M the set of Markov policies $\pi = (f_0, f_1, \dots)$ with $f_n \in F$. Note that ‘Markov’ refers to the fact that the decision at time n depends only on x, μ and z . Further

note that we have $\Pi^M \subset \Pi$ in the following sense: For every $\pi = (f_0, f_1, \dots) \in \Pi^M$ we find a $\sigma = (g_0, g_1, \dots) \in \Pi$ such that

$$\begin{aligned} g_0(x_0) &:= f_0(x_0, \mu_0, 1), \\ g_n(h_n) &:= f_n(x_n, \mu_n(\cdot|h_n), \beta^n), \quad n \in \mathbb{N}. \end{aligned}$$

With this interpretation $V_{n\pi}$ is also defined for $\pi \in \Pi^M$.

Let us now introduce the set

$$\mathcal{C}(E) := \left\{ v : E \rightarrow \mathbb{R} : v \text{ is lower semicontinuous and } v \geq V_0 \right\},$$

where we use the topology of weak convergence on $\mathbb{P}_b(E_Y \times \mathbb{R}_+)$. For $v \in \mathcal{C}(E)$ and $f \in F$ we consider the operator

$$(T_f v)(x, \mu, z) := \int_{E_X} v(x', \Psi(x, f(x, \mu, z), x', \mu, z), \beta z) Q^X(dx'|x, \mu^Y, f(x, \mu, z)), \quad (x, \mu, z) \in E$$

which is well-defined. The minimal cost operator of this Markov Decision Model is given by

$$(Tv)(x, \mu, z) := \inf_{a \in D(x)} \int_{E_X} v(x', \Psi(x, a, x', \mu, z), \beta z) Q^X(dx'|x, \mu^Y, a), \quad (x, \mu, z) \in E \quad (3.7)$$

which is again well-defined and $T_f V_0 \geq TV_0 \geq V_0$ (see also the proof below). Note that $V_0 \in \mathcal{C}(E)$. If a decision rule $f \in F$ is such that $T_f v = Tv$, then f is called a *minimizer* of v . We obtain:

Theorem 3.3. *It holds that*

- a) *For a policy $\pi = (f_0, f_1, f_2, \dots) \in \Pi^M$ we have the following cost iteration:*
 $V_{n\pi} = T_{f_0} \dots T_{f_{n-1}} V_0$ for $n = 1, \dots, N$.
- b) $V_n \in \mathcal{C}(E)$ and $V_n = TV_{n-1}$, for $n = 1, \dots, N$, i.e.,

$$V_{n+1}(x, \mu, z) = \inf_{a \in D(x)} \int_{E_X} V_n(x', \Psi(x, a, x', \mu, z), \beta z) Q^X(dx'|x, \mu^Y, a), \quad (x, \mu, z) \in E.$$

The value function of (3.1) is then given by $J_N(x) = V_N(x, Q_0 \otimes \delta_0, 1)$.

- c) *For every $n = 1, \dots, N$ there exists a minimizer $f_n^* \in F$ of V_{n-1} and $(g_0^*, \dots, g_{N-1}^*)$ with*

$$g_n^*(h_n) := f_{N-n}^*(x_n, \mu_n(\cdot|h_n), \beta^n), \quad n = 0, \dots, N-1$$

is an optimal policy for problem (3.1). Note that the optimal policy consists of decision rules which depend on the current state and the current joint conditional distribution of accumulated cost and hidden state.

Proof. The proof of part a) is by induction. For $n = 1$ we obtain with $a := f_0(x, \mu, z)$:

$$\begin{aligned} T_{f_0} V_0(x, \mu, z) &= \int_{E_X} V_0(x', \Psi(x, a, x', \mu, z), \beta z) Q^X(dx'|x, \mu^Y, a) \\ &= \int_{E_Y} \int_{\mathbb{R}_+} \int_{E_X} \int_{\mathbb{R}_+} U(s') \delta_{s+zc(x, y, a)}(ds') q^X(x'|x, y, a) \lambda(dx') \mu(dy, ds) \\ &= \int_{E_Y} \int_{\mathbb{R}_+} U(s + zc(x, y, a)) \mu(dy, ds) \\ &= V_{1\pi}(x, \mu, z). \end{aligned}$$

Suppose the statement is true for $V_{n\pi}$. In order to ease notation we denote for a policy $\pi = (f_0, f_1, f_2, \dots) \in \Pi^M$ by $\tilde{\pi} = (f_1, f_2, \dots)$ the shifted policy. Moreover let again $a := f_0(x, \mu, z)$.

Then

$$\begin{aligned}
(T_{f_0} \dots T_{f_{n-1}} V_0)(x, \mu, z) &= \int_{E_X} V_{n\pi} \left(x', \Psi(x, a, x', \mu, z), \beta z \right) Q^X(dx' | x, \mu^Y, a) \\
&= \int_{E_X} \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}_{x', y'}^\pi \left[U(s' + z \sum_{k=0}^{n-1} \beta^{k+1} c(X_k, Y_k, A_k)) \right] \Psi(x, a, x', \mu, z) (dy', ds') Q^X(dx' | x, \mu^Y, a) \\
&= \int_{E_X} \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}^\pi \left[U(s' + z \sum_{k=1}^n \beta^k c(X_k, Y_k, A_k)) \middle| X_1 = x', Y_1 = y' \right] \\
&\quad \int_{E_Y} \int_{\mathbb{R}_+} q(x', y' | x, y, a) \delta_{s+zc(x, y, a)}(ds') \mu(dy, ds) \nu(dy') \lambda(dx') \\
&= \int_{E_Y} \int_{E_Y} \int_{E_X} \int_{\mathbb{R}_+} \mathbb{E}^\pi \left[U(s + zc(x, y, a) + z \sum_{k=1}^n \beta^k c(X_k, Y_k, A_k)) \middle| X_1 = x', Y_1 = y' \right] \\
&\quad q(x', y' | x, y, a) \mu(dy, ds) \nu(dy') \lambda(dx') \\
&= \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}_{xy}^\pi \left[U(s + z \sum_{k=0}^n \beta^k c(X_k, Y_k, A_k)) \right] \mu(dy, ds) \\
&= V_{n+1}\pi(x, \mu, z).
\end{aligned}$$

and the statement in part a) is shown.

Next we prove parts b) and c) together. From part a) it follows that for $\pi \in \Pi^M$, the value functions in problem (3.6) indeed coincide with the value functions of the previously defined MDP. From MDP theory it follows in particular that it is enough to consider Markov policies Π^M , i.e., $V_n = \inf_{\sigma \in \Pi} V_{n\sigma} = \inf_{\pi \in \Pi^M} V_{n\pi}$ (see e.g. [13] Theorem 18.4). Next consider functions $v \in \mathcal{C}(E)$. We show that $Tv \in \mathcal{C}(E)$ and that there exists a minimizer for v . Statements b) and c) then follow from Theorem 2.3.8 in [1].

We start by proving that $Q^X(\cdot | x, \mu^Y, a)$ is weakly continuous, i.e., we have to show that

$$(x, \mu, a) \mapsto \int v(x') Q^X(dx' | x, \mu^Y, a) \quad (3.8)$$

is continuous for all $v \in C_b(E_X)$ where $C_b(E_X)$ is the set of bounded, continuous functions on E_X . Obviously $\mu_n \Rightarrow \mu$ implies that $\mu_n^Y \Rightarrow \mu^Y$ where \Rightarrow denotes weak convergence. From our standing assumption (A)(v) it follows that $Q(\cdot | x, y, a)$ is weakly continuous. Hence we obtain from Theorem 17.11 in [13] that the function in (3.8) is continuous.

Next we show that

$$(x, a, x', \mu, z) \mapsto \Psi(x, a, x', \mu, z)$$

is continuous at all points where Ψ is defined, i.e., if $(x_n, a_n, x'_n, \mu_n, z_n)$ converges to (x, a, x', μ, z) in $E_X \times A \times E_X \times \mathbb{P}_b(E_Y \times \mathbb{R}_+) \times (0, 1]$ it follows that $\Psi(x_n, a_n, x'_n, \mu_n, z_n) \Rightarrow \Psi(x, a, x', \mu, z)$ where $(x_n, a_n, x'_n, \mu_n, z_n)$ and (x, a, x', μ, z) are such that $\int_{E_Y} q^X(x'_n | x_n, y, a_n) \mu_n^Y(dy) > 0$ and $\int_{E_Y} q^X(x' | x, y, a) \mu^Y(dy) > 0$. Hence for $v \in C_b(E_Y \times \mathbb{R}_+)$ consider

$$\int_{E_Y} \int_{\mathbb{R}_+} v(y', s') \Psi(x, a, x', \mu, z) (dy', ds').$$

If we plug in the definition of Ψ we get a quotient whose numerator and denominator will be investigated separately. For the numerator we obtain

$$\int_{E_Y} \int_{\mathbb{R}_+} \int_{E_Y} v(y', s + zc(x, y, a)) q(x', y' | x, y, a) \nu(dy') \mu(dy, ds)$$

which is continuous by assumption (A)(iv,v) and Theorem 17.11 in [13]. The denominator

$$\int_{E_Y} q^X(x' | x, y, a) \mu^Y(dy)$$

is continuous in (x, a, x', μ) by the same reasoning. Hence Ψ is continuous.

Now suppose $v \in \mathcal{C}(E)$. Taking into account assumption (A), it obviously follows that $(x, x', a, \mu, z) \mapsto v(x', \Psi(x, a, x', \mu, z), \beta z)$ is lower semicontinuous. Again we apply Theorem 17.11 in [13] to obtain that $(x, \mu, z, a) \mapsto \int v(x', \Psi(x, a, x', \mu, z), \beta z) Q^X(dx'|x, \mu^Y, a)$ is lower semicontinuous. Note here that continuity of Ψ at those points where the denominator is positive is sufficient, since the other points form a Q^X null-set. By Proposition 2.4.3 in [1] it follows that $(x, \mu, z) \mapsto (Tv)(x, \mu, z)$ is lower semicontinuous and there exists a minimizer of v .

The inequality $Tv \geq V_0$ is obtained from

$$\begin{aligned} & \int_{E_X} v(x', \Psi(x, a, x', \mu, z), \beta z) Q^X(dx'|x, \mu^Y, a) \\ & \geq \int_{E_X} \int_{\mathbb{R}_+} U(s') \Psi^Y(x, a, x', \mu, z)(ds') Q^X(dx'|x, \mu^Y, a) \\ & = \int_{E_Y} \int_{\mathbb{R}_+} U(s + zc(x, y, a)) \int_{E_X} q^X(x'|x, y, a) \lambda(dx') \mu(dy, ds) \\ & \geq \int_{E_Y} \int_{\mathbb{R}_+} U(s) \mu(dy, ds) = V_0(x, \mu, z) \end{aligned}$$

which implies the statement. \square

Remark 3.4. Note that $\mu \mapsto V_{n\pi}(x, \mu, z)$ is by definition a linear mapping and thus $\mu \mapsto V_n(x, \mu, z)$ is concave.

Remark 3.5. Since $V_0 \in \mathcal{C}(E)$, $TV_0 \geq V_0$ and since the T -operator is monotone, $V_n = T^n V_0$ is increasing in n .

Remark 3.6. Of course instead of minimizing cost one could also consider the problem of maximizing reward. Suppose that $r : D \rightarrow [\underline{r}, \bar{r}]$ (with $0 < \underline{r} < \bar{r}$) is a one-stage reward function and the problem is

$$J_N(x) := \sup_{\sigma \in \Pi} \int_{E_Y} \mathbb{E}_{xy}^\sigma \left[U \left(\sum_{k=0}^{N-1} r(X_k, A_k) \right) \right] Q_0(dy), \quad x \in E_X. \quad (3.9)$$

It is possible to treat this problem in exactly the same way using straightforward modifications.

4. SOME SPECIAL CASES

4.1. The cost function does not depend on the hidden state. An important special case is obtained when the one-stage cost function does not depend on the hidden state y , i.e., $c(x, y, a) = c(x, a)$. In this case the cost which has accumulated so far is always observable. The recursion for the joint conditional distribution $\mu_n(\cdot|h_n)$ of cost and hidden state simplifies considerable. In order to explain this, we define the operator $\Phi : E_X \times A \times E_X \times \mathbb{P}(E_Y) \rightarrow \mathbb{P}(E_Y)$ by

$$\Phi(x, a, x', \mu)(B) := \frac{\int_B \int_{E_Y} q(x', y'|x, y, a) \mu(dy) \nu(dy')}{\int_{E_Y} q^X(x'|x, y, a) \mu(dy)}, \quad B \in \mathcal{B}(E_Y).$$

Note that Φ is exactly the usual updating (Bayesian) operator which appears in classical POMDP (see e.g. [1], section 5.2). It updates the conditional probability of the unobservable state. In what follows denote by (μ_n^ϕ) the sequence of probability measures on E_Y generated by Φ with $\mu_0^\phi := Q_0$. Then we obtain:

Proposition 4.1. *Suppose $c(x, y, a) = c(x, a)$ is independent of y . Then $\mu_n(\cdot|h_n)$ from (3.3) can be written as*

$$\mu_n(B_1 \times B_2|h_n) = \mu_n^Y(B_1|h_n) \cdot \mu_n^S(B_2|h_n), \quad \text{where } B_1 \times B_2 \in \mathcal{B}(E_Y \times \mathbb{R}_+) \quad (4.1)$$

with $\mu_n^S(\cdot|h_n) = \delta_{\sum_{k=0}^{n-1} \beta^k c(x_k, a_k)}$ and $\mu_n^Y(\cdot|h_n) = \mu_n^\Phi(\cdot|h_n)$.

Proof. The proof is by induction on n . The statement for $n = 0$ is true by definition. Now suppose the statement is true for n . We obtain with $h_{n+1} = (h_n, a_n, x')$, $x_n = x$ and $a_n = a$:

$$\begin{aligned} \mu_{n+1}(B_1 \times B_2 | h_{n+1}) &= \frac{\int_{E_Y} \int_{\mathbb{R}_+} \int_{B_1} \int_{B_2} q(x', y' | x, y, a) \nu(dy') \delta_{s+\beta^n c(x,a)}(ds') \mu_n^Y(dy | h_n) \mu_n^S(ds | h_n)}{\int_{E_Y} q^X(x' | x, y, a) \mu_n^Y(dy | h_n)} \\ &= \frac{\int_{B_1} \int_{E_Y} q(x', y' | x, y, a) \mu_n^Y(dy | h_n) \nu(dy')}{\int_{E_Y} q^X(x' | x, y, a) \mu_n^Y(dy | h_n)} \int_{\mathbb{R}_+} \delta_{s+\beta^n c(x,a)}(B_2) \mu_n^S(ds | h_n) \\ &= \Phi(x, a, x', \mu_n^Y(\cdot | h_n))(B_1) \cdot \delta_{\sum_{k=0}^n \beta^k c(x_k, a_k)}(B_2). \end{aligned}$$

Noting that $\mu_n^Y(\cdot | h_n) = \mu_n^\Phi(\cdot | h_n)$ by the induction hypothesis, the statement follows. \square

Thus, the problem simplifies considerably since instead of probability measures on $\mathcal{B}(E_Y \times \mathbb{R}_+)$ we only need to consider probability measures on $\mathcal{B}(E_Y)$ together with an observable sequence of accumulated cost. We can interpret the embedding MDP as one with state space $E_X \times \mathbb{P}(E_Y) \times \mathbb{R}_+ \times (0, 1]$ and the value iteration reads

$$\begin{aligned} V_0(x, \mu, s, z) &:= U(s) \\ V_{n+1}(x, \mu, s, z) &= \inf_{a \in D(x)} \int V_n(x', \Phi(x, a, x', \mu), s + zc(x, a), \beta z) Q^X(dx' | x, \mu, a), \\ &\quad \text{for } (x, \mu, s, z) \in E_X \times \mathbb{P}(E_Y) \times \mathbb{R}_+ \times (0, 1], \end{aligned}$$

where Φ has been defined in the previous calculation.

Remark 4.2. In case there is no unobservable component, i.e., we have a completely observable risk-sensitive MDP, the updating operator $\Psi : E_X \times A \times E_X \times \mathbb{P}(\mathbb{R}_+) \times (0, 1] \rightarrow \mathbb{P}(\mathbb{R}_+)$ boils down to

$$\Psi(x, a, x', \mu, z)(B) = \int_B \delta_{s+zc(x,a)} \mu(ds), \quad B \in \mathcal{B}(\mathbb{R}_+)$$

and we obtain $\mu_n(B | h_n) = \delta_{\sum_{k=0}^{n-1} \beta^k c(x_k, a_k)}(B)$. Hence the updating process is deterministic and instead of μ we can simply store the accumulated cost so far. The value iteration then reads

$$\begin{aligned} V_0(x, s, z) &= U(s), \quad (x, s, z) \in E_X \times \mathbb{R}_+ \times (0, 1] \\ V_{n+1}(x, s, z) &= \inf_{a \in D(x)} \int V_n(x', s + zc(x, a), z\beta) Q(dx' | x, a), \end{aligned}$$

which is exactly the situation which has been investigated in [2].

4.2. A particular class of partially observable control models. The transition law of the process $(X_n, Y_n)_{n \in \mathbb{N}_0}$ we consider here is quite general. For other general models see Chapter 4 in [12]. All these general models contain in particular the following class which appears very often in applications (in particular this is the starting point in [15, 9]):

$$\begin{aligned} X_{n+1} &= h(Y_n) + \eta_{n+1} \\ Y_{n+1} &= b(Y_n, A_n) + \eta_{n+1} \end{aligned}$$

where (ε_n) is a sequence of independent and identically distributed random variables with density φ_ε and (η_n) is a sequence of independent and identically distributed random variables with density φ_η . Both sequences are assumed to be independent and we assume for simplicity that $E_X = E_Y = \mathbb{R}$. We consider here an additive noise but this can also be part of the functions b and h respectively. The transition law under a policy π is for $B_1, B_2 \in \mathcal{B}(\mathbb{R})$ given by

$$\begin{aligned} Q(B_1 \times B_2 | x, y, a) &= \mathbb{P}(X_{n+1} \in B_1, Y_{n+1} \in B_2 | X_n = x, Y_n = y, A_n = a) \\ &= \mathbb{P}(h(y) + \eta_{n+1} \in B_1, b(y, a) + \varepsilon_{n+1} \in B_2) \\ &= \int_{B_1} \varphi_\eta(w - h(y)) dw \int_{B_2} \varphi_\varepsilon(v - b(y, a)) dv. \end{aligned}$$

According to assumption (A)(v) the resulting density q has to be continuous and bounded in all variables. This is for example satisfied if b, h are continuous and $\varphi_\varepsilon, \varphi_\eta$ are continuous and bounded densities, like e.g. the Gaussian density.

4.3. Total costs criterion. In case $\beta = 1$, the costs are not discounted and we minimize the utility of the total costs

$$\sum_{k=0}^{N-1} c(X_k, Y_k, A_k).$$

In this case the z -component of the iteration in Theorem 3.3 b) does not change. Since in general we start with $z = 1$, we can just skip it and obtain the simpler recursion for $n = 0, \dots, N - 1$

$$\begin{aligned} V_0(x, \mu) &:= \int \int U(s) \mu(dy, ds) \\ V_{n+1}(x, \mu) &= \inf_{a \in D(x)} \int_{E_X} V_n(x', \Psi(x, a, x', \mu)) Q^X(dx' | x, \mu^Y, a), \quad (x, \mu) \in E_X \times \mathbb{P}_b(E_Y \times \mathbb{R}_+), \end{aligned}$$

where $\Psi(x, a, x', \mu) := \Psi(x, a, x', \mu, 1)$ from (3.2). Indeed the z -component is equivalent to the knowledge of the time step but since we would like to consider a general problem it makes sense to introduce this component in the model setup in Section 3.

4.4. Exponential Utility function. In this section we assume now that the utility function has the special form $U(x) = \frac{1}{\gamma} e^{\gamma x}$ with $\gamma \neq 0$. This situation is often referred to as the usual risk-sensitive problem. Partially observable problems in this setting have already been considered in [19, 15, 11, 9, 18, 6]. However still in this case our model is far more general than in the previous literature where the filter is derived with a change of measure technique. As we have shown in (3.3) such a measure transformation is not needed for the computation of the filter.

Our aim is to specialize the value iteration from Theorem 3.3 to this case. In order to do this define for $\mu \in \mathbb{P}_b(E_Y \times \mathbb{R}_+)$:

$$\hat{\mu}(B) := \frac{\int_B \int_{\mathbb{R}_+} e^{\gamma s} \mu(dy, ds)}{\int_{\mathbb{R}_+} e^{\gamma s} \mu^S(ds)}, \quad B \in \mathcal{B}(E_Y) \quad (4.2)$$

which obviously yields a new probability measure on $\mathbb{P}(E_Y)$.

Remark 4.3. From Theorem 3.2 it follows directly that $\hat{\mu}$ has a certain interpretation. We obtain for μ_n from Theorem 3.2 that

$$\int_B \int_{\mathbb{R}_+} e^{\gamma s} \mu_n(dy, ds | h_n) = \mathbb{E}^\pi \left[1_B(Y_n) \cdot e^{\gamma \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k)} \middle| h_n \right].$$

If $\hat{\mu}_n$ is the normalized version of this expression then it coincides with the 'information vector' defined e.g. in [15, 6]. Note that we obtain $\hat{\mu}_n$ in a very natural way as a special case of our general μ_n in Section 3.

Further we can write for $(x, \mu, z) \in E$:

$$\begin{aligned} V_n(x, \mu, z) &= \int_{\mathbb{R}_+} e^{\gamma s} \inf_{\pi} \frac{1}{\gamma} \int_{E_Y} \mathbb{E}_{xy}^\pi \left[\exp \left(\gamma z \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k) \right) \right] \mu(dy, ds) \\ &= \int_{\mathbb{R}_+} e^{\gamma s} \mu^S(ds) \cdot \inf_{\pi} \frac{1}{\gamma} \int_{E_Y} \mathbb{E}_{xy}^\pi \left[\exp \left(\gamma z \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k) \right) \right] \hat{\mu}(dy) \\ &=: \int_{\mathbb{R}_+} e^{\gamma s} \mu^S(ds) \cdot \mathbf{e}_n(x, \hat{\mu}, \gamma z). \end{aligned}$$

Using this representation, the value iteration in Theorem 3.3 can be restricted to the functions \mathbf{e}_n . The state space $E_X \times \mathbb{P}(E_Y) \times (0, 1]$ is much simpler because measures are only concentrated on E_Y .

Theorem 4.4. a) For $(x, \mu, z) \in E_X \times \mathbb{P}(E_Y) \times (0, 1]$ it holds that $\mathbf{e}_0(x, \mu, \gamma z) = \frac{1}{\gamma}$ and for $n = 1, \dots, N$

$$\mathbf{e}_{n+1}(x, \mu, \gamma z) = \inf_{a \in D(x)} \int_{E_X} \mathbf{e}_n(x', \Psi_e(x, a, x', \mu, z), \beta \gamma z) \hat{Q}^X(dx'|x, \mu, a, \gamma z),$$

where for $B_1 \in \mathcal{B}(E_X), B_2 \in \mathcal{B}(E_Y)$

$$\hat{Q}^X(B_1|x, \mu, a, z) := \int_{B_1} \int_{E_Y} e^{zc(x, y, a)} q^X(x'|x, y, a) \mu(dy) \lambda(dx'), \quad (4.3)$$

$$\Psi_e(x, a, x', \mu, z)(B_2) := \frac{\int_{B_2} \int_{E_Y} e^{zc(x, y, a)} q(x', y'|x, y, a) \mu(dy) \nu(dy')}{\int_{E_Y} \int_{E_Y} e^{zc(x, y, a)} q(x', y'|x, y, a) \mu(dy) \nu(dy')}. \quad (4.4)$$

The value function of (3.1) is then given by $J_N(x) = \mathbf{e}_N(x, Q_0, \gamma)$.

b) For every $n = 1, \dots, N$ there exists a minimizer $f_n^* \in F$ of \mathbf{e}_{n-1} and $(g_0^*, \dots, g_{N-1}^*)$ with

$$g_n^*(h_n) := f_{N-n}^*(x_n, \mu_n^e(\cdot|h_n), \gamma \beta^n), \quad n = 0, \dots, N-1$$

is an optimal policy for problem (3.1) where the sequence (μ_n^e) of posterior distributions is generated by the updating operator Ψ_e with $\mu_0^e := Q_0$.

Proof. Let $(x, \mu, z) \in E$. On one hand we have that

$$V_{n+1}(x, \mu, z) = \int_{\mathbb{R}_+} e^{\gamma s} \mu^S(ds) \cdot \mathbf{e}_{n+1}(x, \hat{\mu}, \gamma z),$$

on the other hand we have by Theorem 3.3:

$$\begin{aligned} V_{n+1}(x, \mu, z) &= \inf_{a \in D(x)} \int_{E_X} V_n(x', \Psi(x, a, x', \mu, z), \beta z) Q^X(dx'|x, \mu^Y, a) \\ &= \inf_{a \in D(x)} \int_{E_X} \int_{\mathbb{R}_+} e^{\gamma s'} \Psi^S(x, a, x', \mu, z)(ds') \cdot \mathbf{e}_n(x', \hat{\Psi}(x, a, x', \mu, z), \beta \gamma z) Q^X(dx'|x, \mu^Y, a) \\ &= \inf_{a \in D(x)} \int_{E_X} \int_{E_Y} \int_{\mathbb{R}_+} e^{\gamma s + \gamma z c(x, y, a)} q(x', y'|x, y, a) \mu(dy, ds) \nu(dy') \cdot \\ &\quad \mathbf{e}_n(x', \hat{\Psi}(x, a, x', \mu, z), \beta \gamma z) \lambda(dx') \\ &= \int_{\mathbb{R}_+} e^{\gamma s} \mu^S(ds) \cdot \\ &\quad \inf_{a \in D(x)} \int_{E_X} \int_{E_Y} \int_{\mathbb{R}_+} e^{\gamma z c(x, y, a)} q(x', y'|x, y, a) \hat{\mu}(dy) \nu(dy') \mathbf{e}_n(x', \hat{\Psi}(x, a, x', \mu, z), \beta \gamma z) \lambda(dx') \\ &= \int_{\mathbb{R}_+} e^{\gamma s} \mu^S(ds) \cdot \inf_{a \in D(x)} \int_{E_X} \mathbf{e}_n(x', \hat{\Psi}(x, a, x', \mu, z), \beta \gamma z) \hat{Q}^X(dx'|x, \hat{\mu}, a, \gamma z). \end{aligned}$$

It remains to show that $\hat{\Psi}(x, a, x', \mu, z) = \Psi_e(x, a, x', \hat{\mu}, \gamma z)$ which is defined in (4.4). We obtain for $B \in \mathcal{B}(E_Y)$:

$$\begin{aligned} \hat{\Psi}(x, a, x', \mu, z)(B) &= \frac{\int_B \int_{\mathbb{R}_+} e^{\gamma s'} \Psi(x, a, x', \mu, z)(dy', ds')}{\int_{E_Y} \int_{\mathbb{R}_+} e^{\gamma s'} \Psi(x, a, x', \mu, z)(dy', ds')} \\ &= \frac{\int_B \int_{E_Y} \int_{\mathbb{R}_+} q(x', y'|x, y, a) e^{\gamma s + \gamma z c(x, y, a)} \mu(dy, ds) \nu(dy')}{\int_{E_Y} \int_{E_Y} \int_{\mathbb{R}_+} q(x', y'|x, y, a) e^{\gamma s + \gamma z c(x, y, a)} \mu(dy, ds) \nu(dy')} \\ &= \frac{\int_B \int_{E_Y} q(x', y'|x, y, a) e^{\gamma z c(x, y, a)} \hat{\mu}(dy) \nu(dy')}{\int_{E_Y} \int_{E_Y} q(x', y'|x, y, a) e^{\gamma z c(x, y, a)} \hat{\mu}(dy) \nu(dy')} \\ &= \Psi_e(x, a, x', \hat{\mu}, \gamma z)(B). \end{aligned}$$

Hence part a) is shown. Part b) follows as in Theorem 3.3 c). \square

Remark 4.5. If (μ_n) is generated by Ψ with $\mu_0 := Q_0 \otimes \delta_0$ (note that μ_n are probability measures on $\mathcal{B}(E_Y \times \mathbb{R}_+)$), then $\hat{\mu}(\cdot|h_n) = \mu_n^e(\cdot|h_n)$, i.e., (μ_n^e) is the sequence of information vectors (see Remark (4.3)). The statement follows directly from the proof of the previous theorem.

4.5. Power Utility function. In this section we assume that the utility function has the special form $U(x) = \frac{1}{\gamma}x^\gamma$ with $\gamma \neq 0$. Thus, we obtain:

$$\begin{aligned} V_n(x, \mu, z) &= \inf_{\pi} \frac{1}{\gamma} \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}_{xy}^{\pi} \left[\left(s + z \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k) \right)^\gamma \right] \mu(dy, ds) \\ &= z^\gamma \inf_{\pi} \frac{1}{\gamma} \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}_{xy}^{\pi} \left[\left(\frac{s}{z} + \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k) \right)^\gamma \right] \mu(dy, ds) \\ &= z^\gamma \inf_{\pi} \frac{1}{\gamma} \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}_{xy}^{\pi} \left[\left(\tilde{s} + \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k) \right)^\gamma \right] \tilde{\mu}(dy, d\tilde{s}) \\ &=: z^\gamma d_n(x, \tilde{\mu}), \end{aligned}$$

where $\tilde{\mu}$ is defined by $\tilde{\mu}(B_1 \times B_2) := \mu(B_1 \times \frac{1}{z}B_2)$ for $B_1 \times B_2 \in \mathcal{B}(E_Y \times \mathbb{R}_+)$. Hence $\tilde{\mu} \in \mathbb{P}_b(E_Y \times \mathbb{R}_+)$.

Theorem 4.6. a) For $(x, \mu) \in E_X \times \mathbb{P}_b(E_Y \times \mathbb{R}_+)$ it holds $d_0(x, \mu) := \frac{1}{\gamma} \int \int s^\gamma \mu(dy, ds)$ and for $n = 1, \dots, N$

$$d_{n+1}(x, \mu) = \inf_{a \in D(x)} \beta^\gamma \int_{E_X} d_n(x', \Psi_p(x, a, x', \mu)) Q^X(dx'|x, \mu, a),$$

where for $B \in \mathcal{B}(E_Y \times \mathbb{R}_+)$

$$\Psi_p(x, a, x', \mu)(B) := \frac{\int_{E_Y} \int_{\mathbb{R}_+} \left(\int_B q(x', y'|x, y, a) \nu(dy') \delta_{\frac{s+c(x, y, a)}{\beta}}(ds') \right) \mu(dy, ds)}{\int_{E_Y} q^X(x'|x, y, a) \mu^Y(dy)}.$$

The value function of (3.1) is then given by $J_N(x) = d_N(x, Q_0 \otimes \delta_0)$.

b) For every $n = 1, \dots, N$ there exists a minimizer $f_n^* \in F$ of d_{n-1} and $(g_0^*, \dots, g_{N-1}^*)$ with

$$g_n^*(h_n) := f_{N-n}^*(x_n, \mu_n^p(\cdot|h_n)), \quad n = 0, \dots, N-1$$

is an optimal policy for problem (3.1), where the sequence (μ_n^p) is generated by Ψ_p with $\mu_0^p := Q_0 \otimes \delta_0$.

Proof. On one hand we have shown

$$V_{n+1}(x, \mu, z) = z^\gamma d_{n+1}(x, \tilde{\mu}).$$

On the other hand we obtain with Theorem 3.3

$$\begin{aligned} V_{n+1}(x, \mu, z) &= \inf_{a \in D(x)} \int_{E_X} V_n(x', \Psi(x, a, x', \mu, z), \beta z) Q^X(dx'|x, \mu^Y, a) \\ &= \inf_{a \in D(x)} \beta^\gamma z^\gamma \int_{E_X} d_n(x', \tilde{\Psi}(x, a, x', \mu, z), \beta z) Q^X(dx'|x, \mu^Y, a). \end{aligned}$$

It remains to show that $\tilde{\Psi}(x, a, x', \mu, z) = \Psi_p(x, a, x', \tilde{\mu})$.

Here we obtain for $B \in \mathcal{B}(E_Y \times \mathbb{R}_+)$:

$$\begin{aligned}
\tilde{\Psi}(x, a, x', \mu, z)(B) &= \frac{\int_{E_Y} \int_{\mathbb{R}_+} \left(\int_B q(x', y' | x, y, a) \nu(dy') \delta_{\frac{s}{z} + c(x, y, a)}(ds') \right) \mu(dy, ds)}{\int_{E_Y} \int_{\mathbb{R}_+} \int_{\mathbb{R}_+} \int_{E_Y} q(x', y' | x, y, a) \nu(dy') \delta_{\frac{s}{z} + c(x, y, a)}(ds') \mu(dy, ds)} \\
&= \frac{\int_{E_Y} \int_{\mathbb{R}_+} \left(\int_B q(x', y' | x, y, a) \nu(dy') \delta_{\frac{s}{z} + c(x, y, a)}(ds') \right) \tilde{\mu}(dy, d\tilde{s})}{\int_{E_Y} \int_{\mathbb{R}_+} \int_{\mathbb{R}_+} \int_{E_Y} q(x', y' | x, y, a) \nu(dy') \delta_{\frac{s}{z} + c(x, y, a)}(ds') \tilde{\mu}(dy, d\tilde{s})} \\
&= \Psi_p(x, a, x', \tilde{\mu})(B).
\end{aligned}$$

Hence part a) is shown. Part b) follows as in Theorem 3.3 c). \square

Remark 4.7. If (μ_n) is generated by Ψ with $\mu_0 := Q_0 \otimes \delta_0$, then $\tilde{\mu}(\cdot | h_n) = \mu_n^p(\cdot | h_n)$. The statement follows directly from the proof of the previous theorem.

Remark 4.8. Note that the special case $U(x) = \log(x)$ can be treated similar. It can also be obtained from the power utility case by letting $\gamma \rightarrow 0$.

Remark 4.9. Also the updating operators Ψ_e and Ψ_p simplify considerably if the cost function $c(x, y, a)$ is independent of y (see Section 4.1).

5. APPLICATION: RISK-SENSITIVE BAYESIAN HOUSE SELLING PROBLEM

As an application we consider a risk-sensitive Bayesian extension of the classical house selling problem with finite time horizon. We assume that offers for a house X_0, \dots, X_N arrive independently and are identically distributed with distribution Q_θ . Here $\theta \in \Theta$ is an unknown parameter and Θ is assumed to be a Borel space. Further we assume that Q_θ has a λ -density $q(x|\theta)$ which is continuous in both parameters with compact support. A prior distribution Q_0 for θ is given. As long as offers are rejected an observation cost of $c_\theta > 0$ has to be paid which also depends on θ and cannot be observed. We suppose that c_θ is continuous in θ . When an offer is accepted, the price is obtained and the process ends. If one has not stopped before N , the last offer has to be accepted. The aim is to find the maximal risk-sensitive stopping reward

$$J_N(x) := \sup_{0 \leq \tau \leq N} \int_{\Theta} \mathbb{E}_{x\theta} \left[U(X_\tau - c_\theta \tau) \right] Q_0(d\theta) \quad (5.1)$$

where the supremum is taken over all stopping times τ . Here we assume that $U : \mathbb{R} \rightarrow \mathbb{R}$ is strictly increasing and concave. In order to have a well-defined problem we also assume that $\sup_\theta \mathbb{E}_\theta[X_1^+] < \infty$. This risk-sensitive Bayesian house selling problem can be solved in a similar way as our general model with $Y_n \equiv \theta$ and $E_Y = \Theta$, i.e., the unobservable component is simply the unknown parameter and $c(x, \theta) = c_\theta$ (independent of x). However note that we also have a terminal reward in case we have not stopped before which equals the last offer. Risk-sensitive house selling problems with complete observation have been treated in [16]. A risk-sensitive Bayesian house selling problem has been considered in [4] however with fixed observation costs c (independent of θ). We define the updating operator Ψ for the joint conditional probability of the unknown parameter θ and the accumulated cost so far only in case we do not stop because otherwise the problem ends immediately. Also note that since $\beta = 1$ we can skip the z -component in the state space. Moreover, the i.i.d. assumption on the offers implies that Ψ does not depend on x which is the previous offer. The updating operator is given by

$$\Psi(x', \mu)(B_1 \times B_2) = \frac{\int_{B_1} \int_{\mathbb{R}_-} q(x' | \theta) \delta_{s - c_\theta}(B_2) \mu(d\theta, ds)}{\int_{\Theta} q(x' | \theta) \mu^\Theta(d\theta)}, \quad B_1 \times B_2 \in \mathcal{B}(\Theta \times \mathbb{R}_-).$$

According to Theorem 3.3 we obtain J_N by computing the functions V_n . These are given by

$$\begin{aligned} V_0(x, \mu) &= \int U(x+s) \mu^S(ds) =: U_\mu(x) \\ V_n(x, \mu) &= \max \left\{ U_\mu(x), d_n(\mu) \right\} \end{aligned}$$

with $d_n(\mu) := \int_{\mathbb{R}} V_{n-1}(x', \Psi(x', \mu)) Q^X(dx' | \mu^\Theta)$. We have that $J_N(x) = V_N(x, Q_0 \otimes \delta_0)$. Note that $Q^X(dx' | \mu^\Theta)$ is given by

$$Q^X(B | \mu^\Theta) = \int_B \int_{\Theta} q(x' | \theta) \mu^\Theta(d\theta) \lambda(dx'), \quad B \in \mathcal{B}(E_X).$$

When we define $f_n^*(x, \mu) = \text{stop}$ if $U_\mu(x) \geq d_n(\mu)$ and $(g_0^*, \dots, g_{N-1}^*)$ by

$$g_n^*(h_n) := f_{N-n}^*(x_n, \mu_n(\cdot | h_n)), \quad h_n = (x_0, x_1, \dots, x_n), \quad n = 0, \dots, N-1,$$

then the optimal stopping time for problem (5.1) is given by

$$\tau^* := \inf \{ n \in \mathbb{N}_0 : g_n^*(h_n) = \text{stop} \} \wedge N.$$

Let us now further investigate the optimal stopping time τ^* . As in Section 3 we define by $\mu_n(\cdot | h_n)$ the sequence of conditional probabilities generated by the updating-operator. Then we have

$$g_n^*(h_n) = \text{stop} \quad \Leftrightarrow \quad U_{\mu_n(\cdot | h_n)}(x_n) \geq d_{N-n}(\mu_n(\cdot | h_n)).$$

Since $x \mapsto U_\mu(x)$ is increasing and continuous, the inverse function U_μ^{-1} exists and we obtain

$$g_n^*(h_n) = \text{stop} \quad \Leftrightarrow \quad x_n \geq U_{\mu_n(\cdot | h_n)}^{-1}(d_{N-n}(\mu_n(\cdot | h_n))) =: x_{n,N}^*(\mu_n(\cdot | h_n)).$$

We call $x_{n,N}^*(\cdot)$ *reservation level*. The reservation levels depend on μ_n and U . The optimal stopping time is hence the first time, the offer exceeds the corresponding, history dependent reservation level.

Theorem 5.1. a) *The optimal stopping time for the risk-sensitive Bayesian house selling problem is given by*

$$\tau^* = \inf \{ n \in \mathbb{N}_0 : X_n \geq x_{n,N}^*(\mu_n(\cdot | h_n)) \} \wedge N.$$

b) *The reservation levels can recursively be computed by*

$$\begin{aligned} x_{N-1,N}^*(\mu_{N-1}) &= U_{\mu_{N-1}}^{-1} \circ \int_{\mathbb{R}} \int_{\mathbb{R}} U(x+s) \mu_{N-1}^S(ds) Q^X(dx | \mu_{N-1}^\Theta) \\ x_{n,N}^*(\mu_n) &= U_{\mu_n}^{-1} \circ \int_{\mathbb{R}} U_{\Psi(x, \mu_n)} \left(\max \{ x, x_{n+1,N}^*(\Psi(x, \mu_n)) \} \right) Q^X(dx | \mu_n^\Theta). \end{aligned}$$

Proof. Part a) is clear from the definition and the previous results. Part b) can be shown by inserting the correct definitions. For $n = N-1$ we obtain from the definition of $x_{N-1,N}^*(\mu)$ that

$$x_{N-1,N}^*(\mu) = U_\mu^{-1}(d_1(\mu))$$

with

$$d_1(\mu) = \int V_0(x, \Psi(x, \mu)) Q^X(dx | \mu^\Theta).$$

For $x_{n,N}^*$ we obtain by definition:

$$x_{n,N}^*(\mu) = U_\mu^{-1}(d_{N-n}(\mu)).$$

Further $d_{N-n}(\mu_n)$ can be written as

$$\begin{aligned} d_{N-n}(\mu_n) &= \int_{\mathbb{R}} V_{N-n-1}\left(x, \Psi(x, \mu_n),\right) Q^X(dx|\mu_n^\Theta) \\ &= \int_{\mathbb{R}} \max \left\{ U_{\Psi(x, \mu_n)}(x), d_{N-n-1}(\Psi(x, \mu_n)) \right\} Q^X(dx|\mu_n^\Theta) \\ &= \int_{\mathbb{R}} U_{\Psi(x, \mu_n)}\left(\max \left\{ x, U_{\Psi(x, \mu_n)}^{-1} \circ d_{N-n-1}(\Psi(x, \mu_n)) \right\}\right) Q^X(dx|\mu_n^\Theta) \end{aligned}$$

and the statement follows from the definition of $x_{n,N}^*$. \square

6. INFINITE HORIZON PROBLEMS

Here we consider an infinite time horizon and $\beta \in (0, 1)$, i.e., we are interested in

$$J_\infty(x) := \inf_{\sigma \in \Pi} \int_{E_Y} \mathbb{E}_{xy}^\sigma \left[U \left(\sum_{k=0}^{\infty} \beta^k c(X_k, Y_k, A_k) \right) \right] Q_0(dy), \quad x \in E. \quad (6.1)$$

We will consider concave and convex utility functions separately.

6.1. Concave Utility Function. We first investigate the case of a concave utility function $U : \mathbb{R}_+ \rightarrow \mathbb{R}$. This situation represents a risk seeking decision maker.

In this subsection we use the following notations

$$\begin{aligned} V_{\infty\sigma}(x, \mu, z) &:= \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}_{xy}^\sigma \left[U \left(s + z \sum_{k=0}^{\infty} \beta^k c(X_k, Y_k, A_k) \right) \right] \mu(ds, dy), \\ V_\infty(x, \mu, z) &:= \inf_{\sigma \in \Pi} V_{\infty\sigma}(x, \mu, z), \quad (x, \mu, z) \in E. \end{aligned} \quad (6.2)$$

We are interested in obtaining $V_\infty(x, Q_0 \otimes \delta_0, 1) = J_\infty(x)$. For a stationary policy $\pi = (f, f, \dots) \in \Pi^M$ we write $V_{\infty\pi} = V_f$ and denote

$$\begin{aligned} \bar{b}(\mu, z) &:= \int_{\mathbb{R}_+} U \left(s + \frac{z\bar{c}}{1-\beta} \right) \mu^S(ds), \\ \underline{b}(\mu, z) &:= \int_{\mathbb{R}_+} U \left(s + \frac{z\underline{c}}{1-\beta} \right) \mu^S(ds), \quad (\mu, z) \in \mathbb{P}_b(E_Y \times \mathbb{R}_+) \times [0, 1]. \end{aligned}$$

Then we obtain the main theorem of this section:

Theorem 6.1. *The following statements hold true:*

- a) V_∞ is the unique solution of $v = Tv$ in $\mathcal{C}(E)$ with $\underline{b}(\mu, z) \leq v(x, \mu, z) \leq \bar{b}(\mu, z)$ for T defined in (3.7). Moreover, $T^n V_0 \uparrow V_\infty$, $T^n \underline{b} \uparrow V_\infty$ and $T^n \bar{b} \downarrow V_\infty$ for $n \rightarrow \infty$. The value function of (6.1) is given by $J_\infty(x) = V_\infty(x, Q_0 \otimes \delta_0, 1)$.
- b) There exists a minimizer f^* of V_∞ and (g_0^*, g_1^*, \dots) with

$$g_n^*(h_n) := f^*(x_n, \mu_n(\cdot|h_n), \beta^n)$$

is an optimal policy for (6.1).

Proof. a) We first show that $V_n = T^n V_0 \uparrow V_\infty$ for $n \rightarrow \infty$. To this end note that for $U : \mathbb{R}_+ \rightarrow \mathbb{R}$ increasing and concave we obtain the inequality

$$U(s_1 + s_2) \leq U(s_1) + U'_-(s_1)s_2, \quad s_1, s_2 \geq 0$$

where U'_- is the left-hand side derivative of U which exists since U is concave. Moreover, $U'_-(s) \geq 0$ and U'_- is non-increasing. For $(x, \mu, z) \in E$ and $\sigma \in \Pi$ it holds

$$\begin{aligned}
& V_n(x, \mu, z) \leq V_{n\sigma}(x, \mu, z) \leq V_{\infty\sigma}(x, \mu, z) \\
&= \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}_{xy}^\sigma \left[U \left(s + z \sum_{k=0}^{\infty} \beta^k c(X_k, Y_k, A_k) \right) \right] \mu(dy, ds) \\
&\leq \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}_{xy}^\sigma \left[U \left(s + z \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k) \right) \right] \mu(dy, ds) \\
&\quad + \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}_{xy}^\sigma \left[U'_- \left(s + z \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k) \right) z \sum_{m=n}^{\infty} \beta^k c(X_m, Y_m, A_m) \right] \mu(dy, ds) \\
&\leq V_{n\sigma}(x, \mu, z) + \beta^n \frac{z\bar{c}}{1-\beta} \int_{\mathbb{R}_+} U'_-(s + z\bar{c}) \mu^S(ds) \\
&\leq V_{n\sigma}(x, \mu, z) + \beta^n \frac{z\bar{c}}{1-\beta} U'_-(z\bar{c}) =: V_{n\sigma}(x, \mu, z) + \varepsilon_n(z), \tag{6.3}
\end{aligned}$$

where $\varepsilon_n(z)$ has implicitly been defined in the last equation.

Obviously $\lim_{n \rightarrow \infty} \varepsilon_n(z) = 0$. Taking the infimum over all policies in the preceding inequality yields:

$$V_n(x, \mu, z) \leq V_\infty(x, \mu, z) \leq V_n(x, \mu, z) + \varepsilon_n(z).$$

Letting $n \rightarrow \infty$ yields $V_n = T^n V_0 \uparrow V_\infty$ for $n \rightarrow \infty$. Note that the convergence of $T^n V_0$ is monotone (see Remark 3.5).

By direct inspection we obtain $\underline{b} \leq V_\infty \leq \bar{b}$. We next show that $V_\infty = TV_\infty$. Note that $V_n \leq V_\infty$ for all n . Since T is increasing we have $V_{n+1} = TV_n \leq TV_\infty$ for all n . Letting $n \rightarrow \infty$ implies $V_\infty \leq TV_\infty$. For the reverse inequality recall that $V_n + \varepsilon_n \geq V_\infty$ from (6.3). Applying the T -operator yields $V_{n+1} + \varepsilon_{n+1} = T(V_n + \varepsilon_n) \geq TV_\infty$ and letting $n \rightarrow \infty$ we obtain $V_\infty \geq TV_\infty$. Hence it follows $V_\infty = TV_\infty$.

Next, we obtain

$$\begin{aligned}
(T\bar{b})(\mu, z) &= \inf_{a \in D(x)} \int_{\mathbb{R}_+} U \left(s' + \frac{z\beta\bar{c}}{1-\beta} \right) \Psi^S(x, a, x' \mu, z)(ds') \\
&\leq \int_{\mathbb{R}_+} U \left(s + z\bar{c} + \frac{z\beta\bar{c}}{1-\beta} \right) \mu^S(ds) \\
&= \int_{\mathbb{R}_+} U \left(s + \frac{z\bar{c}}{1-\beta} \right) \mu^S(ds) = \bar{b}(\mu, z).
\end{aligned}$$

Analogously $T\underline{b} \geq \underline{b}$. Thus we get that $T^n \bar{b} \downarrow$ and $T^n \underline{b} \uparrow$ and the limits exist. Moreover, we obtain by iteration:

$$\begin{aligned}
& (T^n \underline{b})(x, \mu, z) = \\
&= \inf_{\pi \in \Pi^M} \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}_{xy}^\pi \left[U \left(s + \frac{z\bar{c}\beta^n}{1-\beta} + z \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k) \right) \right] \mu(dy, ds) \geq (T^n V_0)(x, \mu, z). \\
& (T^n \bar{b})(x, \mu, z) = \\
&= \inf_{\pi \in \Pi^M} \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}_{xy}^\pi \left[U \left(s + \frac{z\bar{c}\beta^n}{1-\beta} + z \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k) \right) \right] \mu(dy, ds).
\end{aligned}$$

Using $U(s_1 + s_2) - U(s_1) \leq U'_-(s_1)s_2$ we obtain:

$$\begin{aligned}
0 &\leq (T^n \bar{b})(x, \mu, z) - (T^n \underline{b})(x, \mu, z) \leq (T^n \bar{b})(x, \mu, z) - (T^n V_0)(x, \mu, z) \\
&\leq \sup_{\pi \in \Pi} \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}_{xy}^\pi \left[U \left(s + \frac{z \bar{c} \beta^n}{1 - \beta} + z \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k) \right) - \right. \\
&\quad \left. U \left(s + z \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k) \right) \right] \mu(dy, ds) \\
&\leq \varepsilon_n(z)
\end{aligned}$$

and the right-hand side converges to zero for $n \rightarrow \infty$. As a result $T^n \bar{b} \downarrow V_\infty$ and $T^n \underline{b} \uparrow V_\infty$ for $n \rightarrow \infty$.

Since V_n is lower semicontinuous, this yields immediately that V_∞ is again lower semicontinuous, thus $V_\infty \in \mathcal{C}(E)$.

For the uniqueness suppose that $v \in \mathcal{C}(E)$ is another solution of $v = Tv$ with $\underline{b} \leq v \leq \bar{b}$. Then $T^n \underline{b} \leq v \leq T^n \bar{b}$ for all $n \in \mathbb{N}$ and since the limit $n \rightarrow \infty$ of the right and left-hand side are equal to V_∞ the statement follows.

- b) The existence of a minimizer follows from our standing assumption (A) as in the proof of Theorem 3.3. From our assumption and the fact that $V_\infty \geq V_0$ we obtain

$$V_\infty = \lim_{n \rightarrow \infty} T_{f^*}^n V_\infty \geq \lim_{n \rightarrow \infty} T_{f^*}^n V_0 = \lim_{n \rightarrow \infty} V_{n(f^*, f^*, \dots)} = V_{f^*} \geq V_\infty$$

where the last equation follows with dominated convergence. Hence (g_0^*, g_1^*, \dots) is optimal for (6.1). \square

6.2. Convex Utility Function. Here we consider the problem with convex utility U . This situation represents a risk averse decision maker. The value functions $V_{n\sigma}, V_n, V_{\infty\sigma}, V_\infty$ are defined as in the previous section.

Theorem 6.2. *Theorem 6.1 also holds for convex U .*

Proof. The proof follows along the same lines as in Theorem 6.1. The only difference is that we have to use another inequality: Note that for $U : \mathbb{R}_+ \rightarrow \mathbb{R}$ increasing and convex we obtain the inequality

$$U(s_1 + s_2) \leq U(s_1) + U'_+(s_1)s_2, \quad s_1, s_2 \geq 0$$

where U'_+ is the right-hand side derivative of U which exists since U is convex. Moreover, $U'_+(s) \geq 0$ and U'_+ is increasing. Thus, we obtain for $(x, \mu, z) \in E$ and $\sigma \in \Pi$:

$$\begin{aligned}
V_n(x, \mu, z) &\leq V_{n\sigma}(x, \mu, z) \leq V_{\infty\sigma}(x, \mu, z) \\
&= \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}_{xy}^\sigma \left[U \left(s + z \sum_{k=0}^{\infty} \beta^k c(X_k, Y_k, A_k) \right) \right] \mu(dy, ds) \\
&\leq \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}_{xy}^\sigma \left[U \left(s + z \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k) \right) \right] + \\
&\quad + \mathbb{E}_{xy}^\sigma \left[U'_+ \left(s + z \sum_{k=0}^{\infty} \beta^k c(X_k, Y_k, A_k) \right) z \sum_{k=n}^{\infty} \beta^k c(X_k, Y_k, A_k) \right] \mu(dy, ds) \\
&\leq V_{n\sigma}(x, \mu, z) + \frac{z \bar{c} \beta^n}{1 - \beta} \int_{\mathbb{R}_+} U'_+ \left(s + \frac{z \bar{c}}{1 - \beta} \right) \mu^S(ds).
\end{aligned}$$

Note that the last inequality follows from the fact that c is bounded from above by \bar{c} . Now denote $\delta_n(\mu, z) := \frac{z \bar{c} \beta^n}{1 - \beta} \int_{\mathbb{R}_+} U'_+ \left(s + \frac{z \bar{c}}{1 - \beta} \right) \mu^S(ds)$. Obviously $\lim_{n \rightarrow \infty} \delta_n(\mu, z) = 0$. Taking the infimum over all policies in the above inequality yields:

$$V_n(x, \mu, z) \leq V_\infty(x, \mu, z) \leq V_n(x, \mu, z) + \delta_n(\mu, z).$$

Letting $n \rightarrow \infty$ yields $T^n V_0 \rightarrow V_\infty$.

Further we have to use the inequality

$$\begin{aligned}
0 &\leq (T^n \bar{b})(x, \mu, z) - (T^n \underline{b})(x, \mu, z) \leq (T^n \bar{b})(x, \mu, z) - (T^n V_0)(x, \mu, z) \\
&\leq \sup_{\pi \in \Pi} \int_{E_Y} \int_{\mathbb{R}_+} \mathbb{E}_x^\pi \left[U \left(s + \frac{z \bar{c} \beta^n}{1 - \beta} + z \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k) \right) - \right. \\
&\quad \left. U \left(s + z \sum_{k=0}^{n-1} \beta^k c(X_k, Y_k, A_k) \right) \right] \mu(dy, ds) \\
&\leq \frac{z \bar{c} \beta^n}{1 - \beta} \int_{\mathbb{R}_+} U'_+ \left(s + \frac{z \bar{c}}{1 - \beta} \right) \mu^S(ds) = \delta_n(\mu, z)
\end{aligned}$$

and the right-hand side converges to zero for $n \rightarrow \infty$. \square

6.3. Exponential Utility. Of course the result for the infinite horizon problem can now be specialized to various situations like in Section 4. This can be done rather straightforward. We only present the case of the exponential utility due to its importance.

Corollary 6.3. *In case $U(x) = \frac{1}{\gamma} e^{\gamma x}$ with $\gamma \neq 0$, we obtain*

- a) $V_\infty(x, \mu, z) = \int e^{\gamma s} \mu^S(ds) \cdot \mathbf{e}_\infty(x, \hat{\mu}, \gamma z)$, $(x, \mu, z) \in E_X \times \mathbb{P}(E_Y \times \mathbb{R}_+) \times (0, 1]$ where $\hat{\mu}$ has been defined in (4.2) and the function \mathbf{e}_∞ is the unique fixed point of

$$\mathbf{e}_\infty(x, \mu, \gamma z) = \inf_{a \in D(x)} \int_{E_X} \mathbf{e}_\infty(x', \Psi_e(x, a, x', \mu, \gamma z), \beta \gamma z) \hat{Q}^X(dx' | x, \mu, a, \gamma z),$$

for $(x, \mu, z) \in E_X \times \mathbb{P}(E_Y) \times (0, 1]$ with $U(\frac{z\bar{c}}{1-\beta}) \leq \mathbf{e}_\infty(x, \mu, \gamma z) \leq U(\frac{z\bar{c}}{1-\beta})$. The value function of (6.1) is then given by $J_\infty(x) = \mathbf{e}_\infty(x, Q_0, \gamma)$.

- b) There exists a minimizer f^* of \mathbf{e}_∞ and (g_0^*, g_1^*, \dots) with

$$g_n^*(h_n) := f^*(x_n, \mu_n^e(\cdot | h_n), \gamma \beta^n)$$

is an optimal policy for (6.1), where the sequence (μ_n^e) of posterior distributions is generated by the updating operator Ψ_e with $\mu_0^e := Q_0$ like in Theorem 4.4.

Acknowledgements: The authors would like to thank three referees for helpful comments and suggestions which improved the presentation of the paper.

REFERENCES

- [1] N. Bäuerle and U. Rieder, Markov Decision Processes with Applications to Finance. Springer-Verlag, Berlin Heidelberg, (2011).
- [2] N. Bäuerle and U. Rieder, More risk-sensitive Markov Decision Processes. Mathematics of Operations Research 39(1), 105-120, (2014).
- [3] N. Bäuerle and A. Jaśkiewicz, Risk-sensitive dividend problems. European Journal of Operational Research 242(1), 161-171, (2015).
- [4] N. Bäuerle and U. Rieder, Partially observable risk-sensitive stopping problems. In: Modern Trends in Controlled Stochastic Processes II (A.B. Piunovskiy ed.) Luniver Press, 12-31, (2015).
- [5] T. Bielecki and S. Pliska, Economic properties of the risk sensitive criterion for portfolio management. Review of Accounting and Finance 2, 3-17, (2003).
- [6] R. Cavazos-Cadena and D. Hernández-Hernández, Successive approximations in partially observable controlled Markov chains with risk-sensitive average criterion. Stochastics 77(6), 537-568, (2005).
- [7] R. Cavazos-Cadena and D. Hernández-Hernández, A Characterization of the Optimal Certainty Equivalent of the Average Cost via the Arrow-Pratt Sensitivity Function. Mathematics of Operations Research 41(1), 224-235, (2016).

- [8] M.H.A. Davis and Sebastien Lleo, Risk-Sensitive Investment Management. World Scientific, (2014).
- [9] Di Masi and L. Stettner, Risk sensitive control of discrete time partially observed Markov processes with infinite horizon. Stochastics 67(3-4), 309-322, (1999).
- [10] W.B. Haskell and R. Jain, A convex analytic approach to risk-aware Markov Decision Processes. SIAM Journal on Control and Optimization 53, 1569-1598, (2015).
- [11] D. Hernández-Hernández, Partially observed control problems with multiplicative cost, In: Stochastic Analysis, Control, Optimization and Applications. Birkhäuser Boston, 41-55, (1999).
- [12] O. Hernández-Lerma, Adaptive Markov control processes. Springer-Verlag, (1989).
- [13] K. Hinderer, Foundations of non-stationary dynamic programming with discrete time parameter. Springer-Verlag, Berlin, (1970).
- [14] R.A. Howard and J.E. Matheson, Risk-sensitive Markov Decision Processes. Management Science 18, 356-369, (1972).
- [15] M.R. James and J.S. Baras and R.J. Elliott, Risk-sensitive control and dynamic games for partially observed discrete-time nonlinear systems. IEEE Transactions on Automatic Control 39(4), 780-792, (1994).
- [16] A. Müller, Expected utility maximization of optimal stopping problems. European Journal of Operational Research 122, 101-114, (2000).
- [17] D. Rhenius, Incomplete information in Markovian decision models. The Annals of Statistics 2, 1327-1334, (1974).
- [18] L. Stettner, Risk sensitive portfolio optimization with completely and partially observed factors. IEEE Transactions on Automatic Control 49(3), 457-464, (2004).
- [19] P. Whittle, Risk-sensitive linear quadratic Gaussian control. Advances in Applied Probability 13, 764-777, (1981).
- [20] A.A. Yushkevich, Reduction of a Controlled Markov Model with Incomplete Data to a Problem with Complete Information in the Case of Borel State and Control Space. Theory of Probability & Its Applications 21(1), 153-158, (1976).

(N. Bäuerle) INSTITUTE FOR STOCHASTICS, KARLSRUHE INSTITUTE OF TECHNOLOGY, D-76128 KARLSRUHE, GERMANY

E-mail address: nicole.baeuerle@kit.edu

(U. Rieder) UNIVERSITY OF ULM, D-89069 ULM, GERMANY

E-mail address: ulrich.rieder@uni-ulm.de